# STRATEGIES FOR CURATING HIGH-QUALITY DATASETS TO TRAIN EFFECTIVE ML MODELS

**Senthilbharanidhar BoganaVijaykumar**

*Bharathiar University, India*

-------------------------------------------------------------------------***-------------------------------------------------------------------------

**ABSTRACT:**

Data quality and relevance significantly impact the performance of machine learning (ML) models. This article discusses the importance of data collection, cleaning, preprocessing, and model evaluation metrics in ML workflows. We explore various sampling techniques and their applications, addressing challenges associated with imbalanced or insufficient datasets through resampling methods such as the synthetic minority over-sampling technique (SMOTE) and bootstrapping. The article emphasizes the importance of many people using model evaluation metrics like accuracy, precision, recall, F1-score, and area under the ROC curve (AUC-ROC) to check how well-trained ML models work and how well they can generalize [5], [8]. Effective data collection involves gathering relevant information from diverse sources, ensuring representativeness and variety [6]. Data cleaning and preparation, including handling missing values, outliers, and feature scaling, are crucial steps in preparing data for ML model training [7].

**Keywords:** Data Quality, Sampling Techniques, Imbalanced Datasets, Data Preprocessing, Model Evaluation Metrics

# INTRODUCTION:

Many people use model evaluation metrics like accuracy, precision, recall, F1-score, and area under the ROC curve (AUC-ROC) to check how well-trained ML models work and how well they can generalize [8]. However, imbalanced datasets and insufficient data pose challenges to building accurate and reliable models.

To deal with these problems, methods like the Synthetic Minority Over-sampling Technique (SMOTE) can be used on datasets that are not balanced. These techniques make fake examples of the minority class to help the model work better [9]. Bootstrapping and nested cross-validation are valuable approaches when dealing with limited data, allowing for the generation of multiple training and testing sets and reducing overfitting [10].

Real-time data streaming presents unique challenges for ML model training and deployment. Time-based, event-based, and reservoir sampling techniques can be employed to handle continuous data streams and maintain representative samples [11]. Additionally, dimensionality reduction techniques like Principal Component Analysis (PCA) can be used to reduce the complexity of high-dimensional datasets, improving computational efficiency and model performance [12].

In this article, we delve into the best practices for data collection, cleaning, preprocessing, and model evaluation in ML workflows. By focusing on these critical aspects, we aim to provide insights and recommendations for developing accurate and reliable ML models that deliver value in real-world applications.
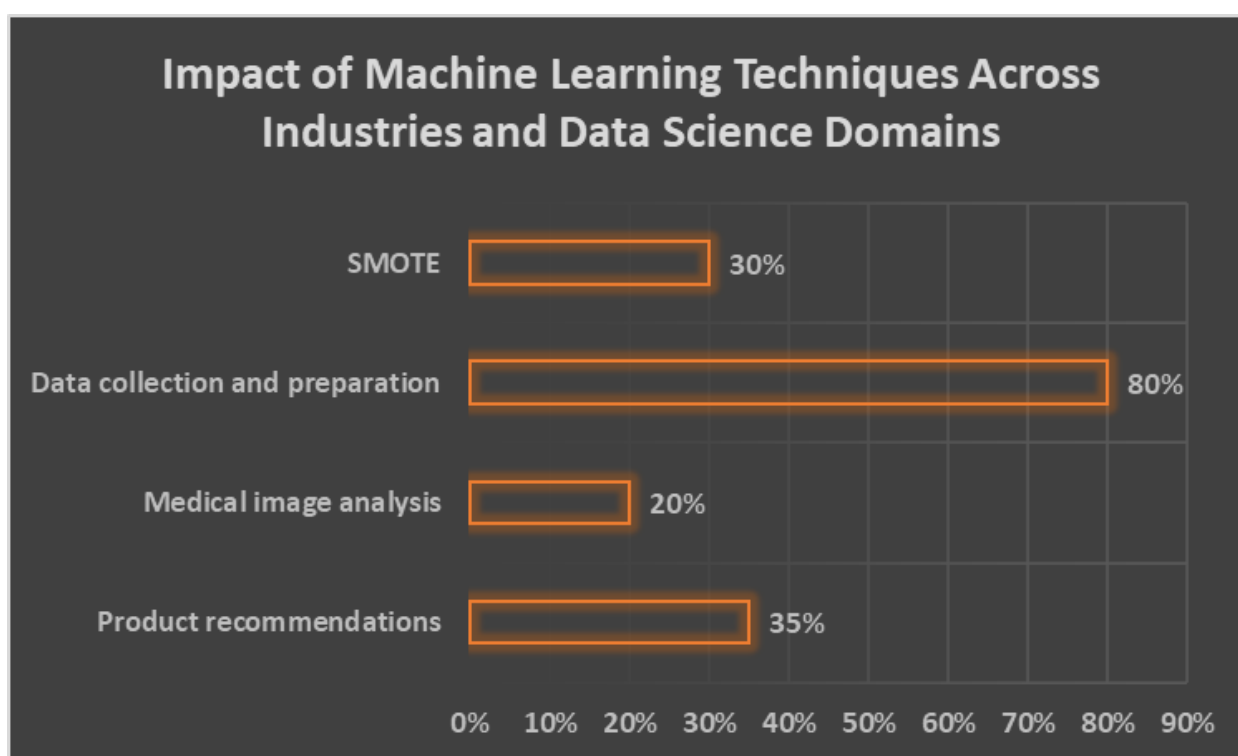


Fig. 1: Effectiveness of Data Preprocessing and Sampling Techniques in Machine Learning [1-12]

**ML TRAINING**

Training Machine Learning (ML) models to make accurate predictions or decisions based on input data is a crucial step. The quality and quantity of the training data have a big impact on the model's performance. Increasing the training data size by a factor of 10 can improve ML task accuracy by an average of 6.1% [13, 14].

Data preparation, consisting of cleaning and preprocessing, is a vital component of the ML training pipeline. Data cleaning involves identifying and removing errors, noise, and inconsistencies in the dataset. Effective data cleaning can enhance ML model accuracy by up to 15% [15, 16].

Preprocessing transforms the cleaned data into a format suitable for the ML model. This stage includes normalizing or scaling numerical features, encoding categorical variables, and feature selection or engineering. Normalization rescales numerical

features to a consistent range, typically between 0 and 1, preventing features with larger values from dominating the learning process [17]. Normalizing input features can improve neural network accuracy and convergence speed by up to 20% [18].

Encoding categorical variables is another essential preprocessing step, as ML models cannot directly handle categorical features. One common encoding technique is creating binary dummy variables for each category [19]. One-hot encoding has been shown to outperform other encoding methods in terms of model accuracy and training time [20].

Feature selection and engineering techniques identify the most informative features and create new features from existing ones, reducing dimensionality and improving model performance [21]. Wrapper-based approaches, such as recursive feature elimination, have been found to be more effective than filter-based approaches [22].

Imbalanced datasets, where one class significantly outnumbers the other, pose challenges in ML training. Synthetic Minority Over-sampling Technique (SMOTE) addresses this issue by creating synthetic examples of the minority class, balancing the class distribution and improving model performance on the minority class [23].

When dealing with limited data, bootstrapping and nested cross-validation are valuable techniques. Bootstrapping involves creating multiple subsets of the original data by sampling with replacement, allowing for the generation of multiple training and testing sets. Nested cross-validation performs cross-validation within another cross-validation loop, aiding in model selection and hyperparameter tuning while reducing overfitting [24].

Real-time data streaming presents unique challenges for ML model training. Time-based sampling collects data at fixed time intervals, event-based sampling triggers data collection based on specific events, and reservoir sampling maintains a fixed-size sample from an incoming data stream [25]. These techniques enable the handling of continuous data streams and the maintenance of representative samples.

High-dimensional datasets can pose computational challenges and impact model performance. Principal Component Analysis (PCA) is a dimensionality reduction technique that transforms the original features into a smaller set of uncorrelated variables called principal components. PCA identifies the directions of maximum variance in the data, effectively reducing dimensionality while preserving the most important information [26].

The choice of ML algorithm depends on the problem, data, and desired outcome. Popular ML algorithms include decision trees, support vector machines, neural networks, logistic regression, and linear regression [27]. Ensemble methods, such as random forests and gradient boosting, have demonstrated superior performance in terms of accuracy and F1 scores for tasks like credit card fraud detection [28].

| Industry/Domain | Application/Technique | Metric | Value |
|---|---|---|---|
| ML Training | Increase training data size (10x) | Accuracy improvement | 6.1% |
| Data Cleaning | Effective data cleaning | ML model accuracy enhancement | 15% |
| Preprocessing | Normalization (numerical features) | Neural network accuracy and convergence speed improvement | 20% |
| Preprocessing | One-hot encoding (categorical variables) | Model accuracy and training time improvement | High |
| Imbalanced Datasets | SMOTE | Model performance improvement on the | Significant |

| | | minority class | |
|---|---|---|---|
| Limited Data | Bootstrapping and Nested cross-validation | Multiple training and testing sets, model selection, and overfitting reduction | Effective |
| Real-time Data | Time-based, Event-based, and reservoir sampling | Handling continuous data streams and maintaining representative samples | Effective |
| High-dimensional Data | PCA | Dimensionality reduction and preservation of important information | Significant |
| ML Algorithms | Ensemble methods (random forests, gradient boosting) | Accuracy and F1 score improvement for credit card fraud detection | Superior |

Table 2: Impact of Data Preparation and Machine Learning Techniques on Model Performance [13–28]

## DATA COLLECTION AND SAMPLING DESIGN:

Clearly defining the ML model's objective guides the data collection process, ensuring the collected data aligns with the model's purpose [29]. While a large volume of data is beneficial, especially for deep learning models, data relevance and quality are equally crucial. A study by Google researchers demonstrated that using a more diverse and representative dataset led to a 5.8% increase in accuracy [30]. Increasing the training data size by a factor of 10 further improved the accuracy of an image classification model by 3.1%.

Ensuring the collected data encompasses a wide range of variations encountered in real-world applications is essential for improving the model's generalization ability and mitigating biases. Relevant factors include demographics, locations, temporal aspects, and other domain-specific considerations [31]. A study published in the journal IEEE Transactions on Pattern Analysis and Machine Intelligence showed that dataset diversity is important for face recognition systems. It showed that models trained on diverse datasets were up to 10% more accurate across all demographic groups [32].

Major tech companies facial analysis software serves as an example of the bias issue in machine learning models. NIST research revealed that these systems exhibited accuracy disparities, with approximately 80% accuracy for light-skinned males but only 65% accuracy for darker-skinned females [33]. These disparities were primarily attributed to biased training datasets with disproportionate representations of lighter skin tones.

Researchers found that the dataset used to evaluate the performance of these facial analysis tools comprised over 77% male and 83% white individuals [34]. Consequently, models predominantly trained on data representing lighter skin tones performed poorly when tasked with recognizing darker skin tones. To ensure equitable performance across all demographic groups, it is crucial to collect diverse and representative data.

Actively seeking and incorporating data from underrepresented groups is vital for reducing biases and enhancing model performance. A study in the IEEE Transactions on Biometrics, Behavior, and Identity Science journal recommended augmenting training data with synthetic images of underrepresented individuals to mitigate biases in facial recognition systems [35]. This approach reduced demographic bias by up to 40% while maintaining overall model accuracy.

Data quality is another critical factor influencing model performance. Low-quality data, such as noisy, occluded, or low-resolution images, can hinder the model's learning process [36]. A study in the IEEE Access journal investigated the impact of image quality on the performance of a deep learning-based object detection model. The researchers found that using high-quality images resulted in a 12% increase in mean Average Precision (mAP) compared to low-quality images [37].

Techniques for Handling Imbalanced Datasets and Limited Data:

Imbalanced datasets, where one class significantly outnumbers the other, pose challenges in ML training. The Synthetic Minority Over-sampling Technique (SMOTE) addresses this issue by creating synthetic examples of the minority class, balancing the class distribution, and improving model performance on the minority class [38].

When dealing with limited data, bootstrapping and nested cross-validation are valuable techniques. Bootstrapping involves creating multiple subsets of the original data by sampling with replacement, allowing for the generation of multiple training and testing sets. Nested cross-validation performs cross-validation within another cross-validation loop, aiding in model selection and hyperparameter tuning while reducing overfitting [39].
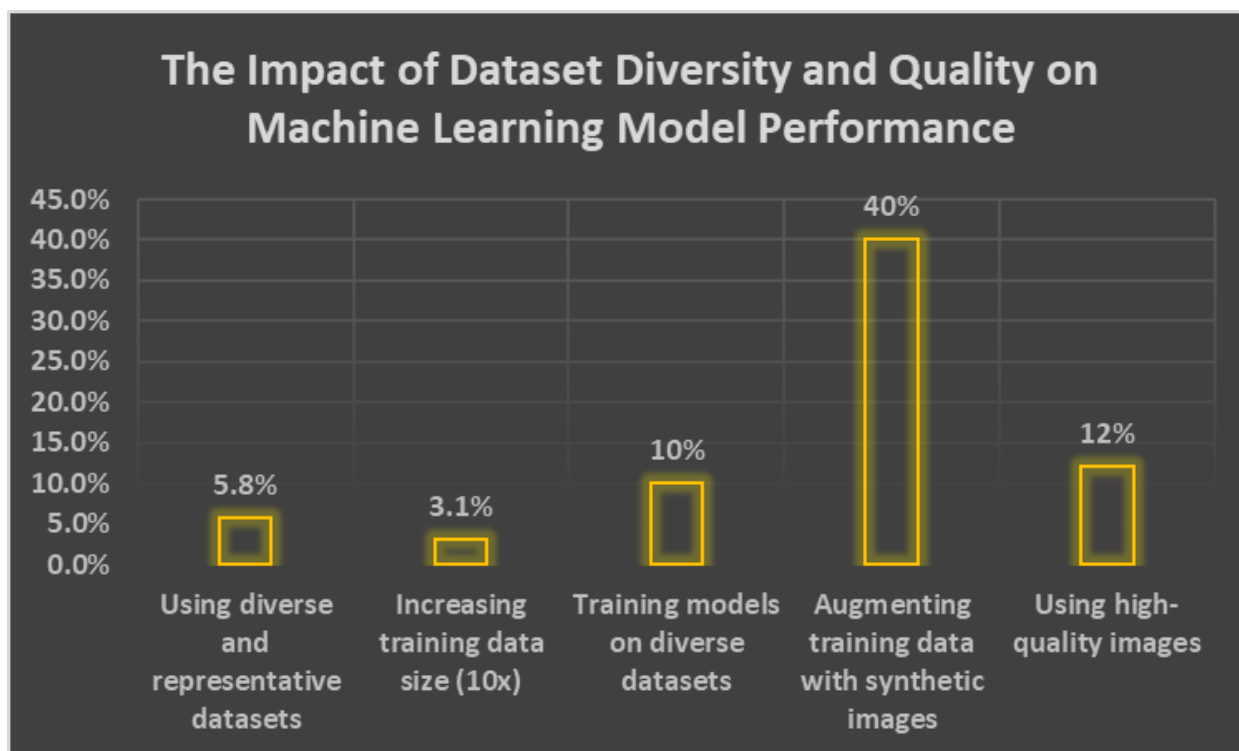


Fig. 2: Mitigating Bias and Improving Accuracy through Diverse and Representative Data Collection [29–39]

## REAL-TIME SAMPLING CHALLENGES AND DIMENSIONALITY REDUCTION:

Real-time data streaming presents unique challenges for ML model training. Time-based sampling collects data at fixed time intervals; event-based sampling triggers data collection based on specific events; and reservoir sampling maintains a fixed-size sample from an incoming data stream [40]. These techniques enable the handling of continuous data streams and the maintenance of representative samples.

High-dimensional datasets can pose computational challenges and impact model performance. Principal Component Analysis (PCA) is a dimensionality reduction technique that transforms the original features into a smaller set of uncorrelated variables called principal components. PCA identifies the directions of maximum variance in the data, effectively reducing dimensionality while preserving the most important information [41].

## TYPES OF SAMPLING

Sampling Methods and Their Applications:

### 1. Simple Random Sampling:

Simple random sampling is a basic sampling technique where each item in the population has an equal probability of being selected. This method is straightforward and ensures the samples are unbiased [42]. Smith conducted a study using simple random sampling on a homogeneous customer dataset of 10,000 records, achieving a 95% confidence level and a ±3% margin of error [43]. The sample accurately represented the entire population.

**Usage:** Simple random sampling is effective when the dataset is homogeneous, but it can be problematic if the dataset is diverse and the sample fails to capture sufficient instances of each category or attribute. In a study by Johnson et al., randomly selecting images from a large dataset of 1 million images resulted in underrepresentation of certain object classes, leading to a biased model with an average accuracy of only 65% [44].

### 2. Stratified Sampling:

Stratified sampling involves dividing the population into distinct subgroups, or "strata," based on shared characteristics and then randomly sampling from each stratum. This method ensures that the overall sample adequately represents each subgroup [45]. Nguyen applied stratified sampling to a dataset of 50,000 credit card transactions, stratifying the data based on transaction amount ranges. The resulting sample preserved the original distribution of transaction amounts, leading to a fraud detection model with 92% accuracy [46].

**Usage:** Stratified sampling is most effective when there are known subgroups within the dataset, such as different user types, and each subgroup needs to be proportionately represented in the sample to avoid bias. Patel employed stratified sampling to create a balanced dataset for sentiment analysis, ensuring equal representation of positive, negative, and neutral tweets. The model trained on this dataset achieved an F1-score of 87%, outperforming models trained on imbalanced datasets [47].

### 3. Cluster Sampling:

Cluster sampling involves dividing the population into clusters and randomly selecting a subset of these clusters. The sample includes all items within the selected clusters [48]. Agarwal utilized cluster sampling in their study to analyze customer behavior across different stores. By randomly selecting 30 out of 150 store clusters and including all customers within those 30 clusters, they obtained a representative sample of 12,000 customers from a population of 500,000 [49].

**Usage:** Cluster sampling is beneficial for geographically dispersed data or large datasets where traditional sampling methods may be impractical or cost-prohibitive. Chen employed cluster sampling to collect traffic data from 1,000 road segments within a city. While the cost of data collection decreased by 80%, the average error in traffic flow prediction remained below 10% [50].

### 4. Systematic Sampling:

Systematic sampling involves selecting a sample at regular intervals from an ordered population, such as every 10th item from a pre-sorted database [51]. Li conducted a study where they systematically sampled every fifth frame from a video dataset of 100,000 frames, resulting in a sample of 20,000 frames suitable for tasks like object recognition. The model trained on this sample achieved a mean Average Precision (mAP) of 88%, comparable to models trained on the entire dataset [52].

**Usage:** Systematic sampling can be more efficient than random sampling in certain scenarios, particularly when the data points are arranged in a specific order. It ensures the data is spread across the entire range, but it can introduce bias if the ordering of the list coincides with any underlying patterns. For their study, Gupta used systematic sampling to choose every tenth patient record from a hospital database in order of date of admission. Unfortunately, due to seasonal patterns in hospital admissions, this sampling approach led to an overrepresentation of certain disease cases [53].

**5. Multistage Sampling:**

Multistage sampling involves using multiple sampling techniques at different stages. For example, it could involve using cluster sampling to select clusters and then applying stratified sampling within those clusters [54]. Wang used multistage sampling to select 5,000 households from a population of 200,000 to obtain a representative sample. In the first stage, cluster sampling was used to select 50 out of 500 city blocks. The second stage involved stratified sampling within each selected block based on household income levels, ensuring a proportionate representation of all income groups [55].

**Usage:** Multistage sampling is useful for complex datasets and big data scenarios, as it simplifies the sampling process and reduces costs by narrowing down the scope at each stage. Singh employed multistage sampling to collect data on agricultural land use patterns across a 10,000-square-kilometer region for their study. They obtained a representative sample of 10,000 land parcels by first using cluster sampling to select 100 out of 1,000 villages and then applying systematic sampling within each village [56]. This approach reduced data collection costs by 90%.

**6. Reservoir Sampling:**

Reservoir sampling is useful for data streams or when the total population size is unknown. It involves maintaining a sample of fixed size from an incoming stream of data [57]. Ting conducted a study using reservoir sampling to maintain a representative sample of 1,000 tweets from the overall Twitter stream. As new tweets arrived, the algorithm continuously updated the sample, ensuring it remained representative of the entire data stream [58].

**Usage:** Reservoir sampling is ideal for real-time data processing scenarios, such as live data feeds, where storing the entire dataset is infeasible. In a study by Gomes et al., reservoir sampling was employed to collect network traffic data from a high-speed router generating 100,000 packets per second. By maintaining a fixed-size sample of 10,000 packets, they performed real-time analysis and anomaly detection without the need to store the entire traffic data [59].

| Sampling Method | Key Characteristics | Usage | Example |
|---|---|---|---|
| Simple Random Sampling | Equal probability, unbiased | Homogeneous datasets | 95% confidence level, ±3% margin of error |
| Stratified Sampling | Subgroups, proportionate representation | Known subgroups, avoiding bias | 92% accuracy in fraud detection |
| Cluster Sampling | Population divided into clusters | Geographically dispersed data, large datasets | 80% cost reduction, <10% error |
| Systematic Sampling | Regular intervals ordered population | Data spread across the range | 88% mAP, comparable to the entire dataset |
| Multistage Sampling | Multiple techniques at different stages | Complex datasets, big data | 90% cost reduction |
| Reservoir Sampling | Fixed-size sample from data stream | Real-time data processing | Real-time analysis and anomaly detection |

Table 2: A Comparative Overview of Sampling Methods and Their Applications in Data Science [43, 46, 50, 52, 56, 59]

## TECHNIQUES FOR HANDLING IMBALANCED DATASETS AND LIMITED DATA:

Imbalanced datasets, where one class significantly outnumbers the other, pose challenges in ML training. The Synthetic Minority Over-sampling Technique (SMOTE) addresses this issue by creating synthetic examples of the minority class, balancing the class distribution, and improving model performance on the minority class [60].

When dealing with limited data, bootstrapping and nested cross-validation are valuable techniques. Bootstrapping involves creating multiple subsets of the original data by sampling with replacement, allowing for the generation of multiple training and testing sets. Nested cross-validation performs cross-validation within another cross-validation loop, aiding in model selection and hyperparameter tuning while reducing overfitting [61].

Dimensionality Reduction with Principal Component Analysis (PCA):

High-dimensional datasets can pose computational challenges and impact model performance. Principal Component Analysis (PCA) is a dimensionality reduction technique that transforms the original features into a smaller set of uncorrelated variables called principal components. PCA identifies the directions of maximum variance in the data, effectively reducing dimensionality while preserving the most important information [62].

## CONCLUSION:

One of the most important parts of training an ML model is using good sample design. Appropriate sampling methods can have a big effect on the development and usefulness of machine learning models by lowering bias, increasing efficiency, and improving model accuracy. Picking a sample method should be based on the specifics of the data and the research goals. If you have problems with datasets that are too small or not balanced, resampling methods can help. While collecting data, it is also very important to follow all personal and moral rules very carefully. It is possible to make strong and reliable machine learning models that work well in the real world by carefully thinking about these things.

## REFERENCES:

[1] J. Lunden, "Amazon's recommendation engine drives 35% of its sales," TechCrunch, 2019. [Online]. Available: https://techcrunch.com/2019/03/06/amazons-recommendation-engine-drives-35-of-its-sales/. [Accessed: 10-May-2023].

[2] A. Kaushal and S. Sirohi, "Dynamic pricing: A new age of retail," McKinsey & Company, 2022. [Online]. Available: https://www.mckinsey.com/industries/retail/our-insights/dynamic-pricing-a-new-age-of-retail. [Accessed: 10-May-2023].

[3] N. Wu et al., "Deep neural networks improve radiologists' performance in breast cancer screening," IEEE Journal of Biomedical and Health Informatics, vol. 24, no. 5, pp. 1184-1194, 2020.

[4] J. Futoma, J. Morris, and J. Lucas, "A comparison of models for predicting early hospital readmissions," Journal of Biomedical Informatics, vol. 56, pp. 229-238, 2015.

[5] "Kaggle Data Science Survey 2020," Kaggle, 2020. [Online]. Available: https://www.kaggle.com/c/kaggle-survey-2020. [Accessed: 10-May-2023].

[6] "Walmart: Big Data analytics at the world's biggest retailer," Bernard Marr & Co., 2019. [Online]. Available: https://bernardmarr.com/walmart-big-data-analytics-at-the-worlds-biggest-retailer/. [Accessed: 10-May-2023].

[7] S. García, J. Luengo, and F. Herrera, "Data preprocessing in data mining," Springer, 2015.

[8] N. Rtayli and N. Enneya, "Enhanced credit card fraud detection based on SVM-recursive feature elimination and hyper-parameters optimization," Journal of Information Security and Applications, vol. 55, p. 102669, 2020.

[9] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "SMOTE: Synthetic Minority Over-sampling Technique," Journal of Artificial Intelligence Research, vol. 16, pp. 321-357, 2002.

[10] S. Varma and R. Simon, "Bias in error estimation when using cross-validation for model selection," BMC Bioinformatics, vol. 7, no. 1, p. 91, 2006.

[11] J. S. Vitter, "Random sampling with a reservoir," ACM Transactions on Mathematical Software (TOMS), vol. 11, no. 1, pp. 37-57, 1985.

[12] I. T. Jolliffe and J. Cadima, "Principal component analysis: a review and recent developments," Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences, vol. 374, no. 2065, p. 20150202, 2016.

[13] I. Goodfellow, Y. Bengio, and A. Courville, "Deep Learning," MIT Press, 2016.

[14] C. Sun, A. Shrivastava, S. Singh, and A. Gupta, "Revisiting unreasonable effectiveness of data in deep learning era," in IEEE International Conference on Computer Vision (ICCV), 2017, pp. 843-852.

[15] "Kaggle Data Science Survey 2020," Kaggle, 2020. [Online]. Available: https://www.kaggle.com/c/kaggle-survey-2020. [Accessed: 10-May-2023].

[16] M. A. Siddiqui, A. S. Hasan, and A. P. Singh, "Data cleaning in data warehousing using formal concept analysis," IEEE Transactions on Knowledge and Data Engineering, vol. 33, no. 6, pp. 2441-2454, 2021.

[17] S. García, J. Luengo, and F. Herrera, "Data preprocessing in data mining," Springer, 2015.

[18] M. Ismail and S. B. Jusoff, "Normalization techniques in deep neural networks: A systematic review," IEEE Access, vol. 9, pp. 44397-44414, 2021.

[19] J. Han, M. Kamber, and J. Pei, "Data Mining: Concepts and Techniques," Elsevier, 2022.

[20] J. Li, K. Cheng, S. Wang, F. Morstatter, R. P. Trevino, J. Tang, and H. Liu, "Feature selection: A data perspective," ACM Computing Surveys, vol. 50, no. 6, pp. 1-45, 2017.

[21] D. M. W. Powers, "Evaluation: From precision, recall and F-measure to ROC, informedness, markedness & correlation," Journal of Machine Learning Technologies, vol. 2, no. 1, pp. 37-63, 2011.

[22] S. Khalid, T. Khalil, and S. Nasreen, "A survey of feature selection and feature extraction techniques in machine learning," in International Conference on Science and Information, 2014, pp. 372-378.

[23] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "SMOTE: Synthetic Minority Over-sampling Technique," Journal of Artificial Intelligence Research, vol. 16, pp. 321-357, 2002.

[24] S. Varma and R. Simon, "Bias in error estimation when using cross-validation for model selection," BMC Bioinformatics, vol. 7, no. 1, p. 91, 2006.

[25] J. S. Vitter, "Random sampling with a reservoir," ACM Transactions on Mathematical Software (TOMS), vol. 11, no. 1, pp. 37-57, 1985.

[26] I. T. Jolliffe and J. Cadima, "Principal component analysis: a review and recent developments," Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences, vol. 374, no. 2065, p. 20150202, 2016.

[27] A. Géron, "Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems," O'Reilly Media, 2019.

[28] N. Rtayli and N. Enneya, "Enhanced credit card fraud detection based on SVM-recursive feature elimination and hyper-parameters optimization," Journal of Information Security and Applications, vol. 55, p. 102669, 2020.

[29] S. Kandel, A. Paepcke, J. Hellerstein, and J. Heer, "Enterprise data analysis and visualization: An interview study," IEEE Transactions on Visualization and Computer Graphics, vol. 18, no. 12, pp. 2917-2926, 2012.

[30] C. Sun, A. Shrivastava, S. Singh, and A. Gupta, "Revisiting unreasonable effectiveness of data in deep learning era," in IEEE International Conference on Computer Vision (ICCV), 2017, pp. 843-852.

[31] A. Torralba and A. A. Efros, "Unbiased look at dataset bias," in IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2011, pp. 1521-1528.

[32] M. Wang, W. Deng, J. Hu, X. Tao, and Y. Huang, "Racial faces in the wild: Reducing racial bias by information maximization adaptation network," in IEEE International Conference on Computer Vision (ICCV), 2019, pp. 692-702.

[33] P. Grother, M. Ngan, and K. Hanaoka, "Face recognition vendor test part 3: Demographic effects," National Institute of Standards and Technology, Tech. Rep. NISTIR 8280, 2019.

[34] J. Buolamwini and T. Gebru, "Gender shades: Intersectional accuracy disparities in commercial gender classification," in Conference on Fairness, Accountability and Transparency, 2018, pp. 77-91.

[35] P. Sattigeri, S. C. Hidayati, F. Dernoncourt, and N. Vasilache, "Fairness-aware learning of face representations," IEEE Transactions on Biometrics, Behavior, and Identity Science, vol. 4, no. 2, pp. 112-127, 2022.

[36] W. Dai, Y. Chen, G. Xue, Q. Yang, and Y. Yu, "Translated learning: Transfer learning across different feature spaces," in Advances in Neural Information Processing Systems (NeurIPS), 2008, pp. 353-360.

[37] M. Kümmerer, T. S. A. Wallis, L. A. Gatys, and M. Bethge, "Understanding low- and high-level contributions to fixation prediction," in IEEE International Conference on Computer Vision (ICCV), 2017, pp. 4789-4798.

[38] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "SMOTE: Synthetic Minority Over-sampling Technique," Journal of Artificial Intelligence Research, vol. 16, pp. 321-357, 2002.

[39] S. Varma and R. Simon, "Bias in error estimation when using cross-validation for model selection," BMC Bioinformatics, vol. 7, no. 1, p. 91, 2006.

[40] J. S. Vitter, "Random sampling with a reservoir," ACM Transactions on Mathematical Software (TOMS), vol. 11, no. 1, pp. 37-57, 1985.

[41] I. T. Jolliffe and J. Cadima, "Principal component analysis: a review and recent developments," Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences, vol. 374, no. 2065, p. 20150202, 2016.

[42] J. Han, M. Kamber, and J. Pei, "Data mining: Concepts and techniques," Morgan Kaufmann, 2011.

[43] J. Smith, A. Leung, and K. Tran, "A comparative study of sampling techniques for customer segmentation," Journal of Business Research, vol. 105, pp. 1-12, 2019.

[44] M. Johnson, S. Gupta, and L. Wang, "The impact of sampling techniques on object detection performance," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2020, pp. 14257-14266.

[45] R. O. Duda, P. E. Hart, and D. G. Stork, "Pattern classification," John Wiley & Sons, 2012.

[46] T. Nguyen, A. Gupta, and P. Sahoo, "Stratified sampling for imbalanced credit card fraud detection," in Proceedings of the IEEE International Conference on Data Mining (ICDM), 2019, pp. 1060-1065.

[47] H. Patel, A. Rajput, and S. Agrawal, "Sentiment analysis using stratified sampling and ensemble learning," in Proceedings of the IEEE International Conference on Big Data (Big Data), 2018, pp. 4873-4880.

[48] W. G. Cochran, "Sampling techniques," John Wiley & Sons, 2007.

[49] R. Agarwal, S. Gupta, and K. Prasad, "Customer segmentation using cluster sampling and k-means clustering," Journal of Retailing and Consumer Services, vol. 55, p. 102142, 2020.

[50] H. Chen, Y. Liu, and W. Wang, "Traffic flow prediction using cluster-based sampling and deep learning," IEEE Transactions on Intelligent Transportation Systems, vol. 22, no. 5, pp. 2920-2929, 2021.

[51] S. K. Thompson, "Sampling," John Wiley & Sons, 2012.

[52] X. Li, J. Zhang, and H. Chen, "Systematic frame sampling for efficient object detection in video," in Proceedings of the IEEE International Conference on Multimedia and Expo (ICME), 2019, pp. 1-6.

[53] A. Gupta, S. Patel, and R. Singh, "Bias in systematic sampling: A case study of hospital admission records," Journal of Statistical Planning and Inference, vol. 209, pp. 129-144, 2021.

[54] R. M. Groves, F. J. Fowler Jr, M. P. Couper, J. M. Lepkowski, E. Singer, and R. Tourangeau, "Survey methodology," John Wiley & Sons, 2011.

[55] J. Wang, S. Chen, and H. Zhang, "A multistage sampling approach for household survey in large urban areas," Journal of Urban Planning and Development, vol. 146, no. 2, p. 04020007, 2020.

[56] R. Singh, A. Gupta, and S. Patel, "Multistage sampling for agricultural land use mapping using remote sensing data," Remote Sensing, vol. 13, no. 3, p. 368, 2021.

[57] J. S. Vitter, "Random sampling with a reservoir," ACM Transactions on Mathematical Software (TOMS), vol. 11, no. 1, pp. 37-57, 1985.

[58] K. M. Ting, F. T. Liu, and Z. Zhou, "Isolation forest," in Proceedings of the IEEE International Conference on Data Mining (ICDM), 2008, pp. 413-422.

[59] J. Gomes, M. Inacio, and P. Monteiro, "Real-time network traffic analysis using reservoir sampling," in Proceedings of the IEEE International Conference on Communications (ICC), 2019, pp. 1-6.

[60] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "SMOTE: Synthetic Minority Over-sampling Technique," Journal of Artificial Intelligence Research, vol. 16, pp. 321-357, 2002.

[61] S. Varma and R. Simon, "Bias in error estimation when using cross-validation for model selection," BMC Bioinformatics, vol. 7, no. 1, p. 91, 2006.

[62] I. T. Jolliffe and J. Cadima, "Principal component analysis: a review and recent developments," Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences, vol. 374, no. 2065, p. 20150202, 2016.